

## Abstract

Ranking a set of objects involves establishing an order allowing for comparisons between any pair of objects in the set. Oftentimes, due to the unavailability of a ground truth of ranked orders, researchers resort to obtaining judgments from multiple annotators followed by inferring the ground truth based on the collective knowledge of the crowd. However, the aggregation is often ad-hoc and involves imposing stringent assumptions in inferring the ground truth (e.g. majority vote). In this work, we propose Expectation-Maximization (EM) based algorithms that rely on the judgments from multiple annotators and the object attributes for inferring the latent ground truth. The algorithm learns the relation between the latent ground truth and object attributes as well as annotator specific “probabilities of flipping”, a metric to assess annotator quality. We further extend the EM algorithm to allow for a variable “probability of flipping” based on the pair of objects at hand. We test our algorithms on two data sets with synthetic annotations and investigate the impact of annotator quality and quantity on the inferred ground truth. We also obtain the results on two other data sets with annotations from machine/human annotators and interpret the output trends based on the data characteristics.

***Index terms***— Learning to Rank, Expectation Maximization, Multiple Annotators, Support Vector Ranker

# Inferring object rankings based on noisy pairwise comparisons from multiple annotators

Rahul Gupta, Shrikanth Narayanan

December 19, 2016

## 1 Introduction

Given a set of items, ranking involves establishing a partial order over the items. This ordering allows comparison between two items, in which the first is either ranked higher, lower or equal to the second [1]. This is commonly termed as a pairwise approach and has been investigated in relation to information retrieval [2], ranking web pages [3] and even analysis of human behavioral constructs such as emotions [4]. Within the problem of modeling preferences using pairwise comparisons, inferring the true order given comparisons from noisy annotators [5] is very relevant. Often, due to the unavailability of the ground truth, experimenters resort to accumulating judgments from multiple annotators and performing a fusion of their collective knowledge. This trend has existed beyond learning to rank and has also been observed in classification and regression tasks [6]. Particularly within the domain of classification, several researchers have proposed novel ways of jointly modeling the annotators in inferring the latent ground truth [7, 8]. Although prior research has addressed similar problems within ranking, the methods enforce a specific structure (e.g., Borda count method, Nanson method [9, 10]) on annotator judgments in inferring the latent ground truth. In this work, we present Expectation-Maximization (EM) [11] based algorithms inspired from work in classification problems to infer the latent ground truth in ranking objects. Through these algorithms, we not only aim to relax the ad-hoc constraints imposed in ground truth computation of preferences but also open up possibilities to integrate the existing approaches within ranking and classification addressing similar problems.

Given noisy pairwise preferences from multiple annotators, the proposed algorithms target to infer a single ground truth ranking while also computing a reliability metric for each of the annotators. We assume the ground truth to be a latent variable that can be inferred not only based on the noisy

pairwise comparisons from multiple annotators, but also the distribution of a set of attributes/features corresponding to the pair of items being compared. We approach this problem using the Expectation-Maximization (EM) framework [12] and develop a Joint Annotator Modeling (JAM) scheme, inspired from existing literature in modeling multiple annotators [7, 8]. The JAM schemes assume that, given the set of attributes/features for a pair of objects, there exists a latent true preference order. Furthermore, the annotators either retain or flip this preference order based on an annotator-specific reliability metric, the “probability of flipping”. The JAM scheme initially learns the relationship between the attributes of the object pair and the latent ground truth as well as each annotator’s “probability of flipping”. The final inference on the preference ground truth is made jointly taking into account the model’s belief based on the object attributes and the annotators’ preferences. We further modify the JAM scheme to allow for non-constant “probability of flipping” based on the pair of objects at hand, termed as Variable Reliability Joint Annotator Modeling (VRJAM) scheme. We compare the JAM and VRJAM schemes to existing methods such as majority voting and fusion after Independent Annotator Modeling (IAM) (similar to Borda count method [9]). We evaluate our models on two data sets with synthetic annotations to investigate the impact of annotator quality and quantity on our models. We also evaluate our models on two other data sets with annotations from machines (ground truths available) and humans (ground truths not available). We interpret the outcomes of the models based on the data characteristics and suggest a few future directions. In the next section, we provide a background of the relevant work, followed by the description of various methodologies for inferring latent true preference order from noisy annotator preferences.

## 2 Previous work

Several researchers have addressed the problem of learning to rank from pairwise comparisons with applications to a variety of domains. In particular, works by Hüllermeier and Fürnkranz et al. [1, 13] provide a comprehensive background on preference learning using the pairwise approach. Considering consolidation of other machine learning topics within the framework of ranking, Brinker et al. [14] and Long et al. [15] integrated active learning in ranking problems, Chu et al. [16] provided an extension of Gaussian processes for ranking and He et al. [17] used manifold based ranking for image retrieval. Other notable works proposing novel methods and applications for ranking include learning to rank using non-smooth cost functions [18], the Mcrank algorithm [19] and learning to rank with partially labeled data [20]. Whereas several existing works have addressed other interesting flavors of learning to rank [21], rank aggregation [22] is possibly one of the most well

studied fields under this domain. A prominent setting under rank aggregation is learning a probability distribution centered around a single or a mixture of global rankings. Several works [23–25] present algorithms for rank aggregation using non-negative matrix factorization, nuclear norm minimization and sparse decomposition techniques.

A different problem setting under learning to rank is inferring a ground truth ranking from a set of pairwise preferences available from multiple annotators. Chen et al. [26] address this problem and present an active learning framework that selects a pair of objects as well as the annotator to be queried while training a ranking model. Along similar lines, Kumar et al. [27] investigated algorithms to fuse ranking models trained using noisy crowd. The formulation of inferring latent ground truth from noisy annotations is particularly well studied in classification and regression problems. Dawid et al. [28] presented one of the earlier works in fusion annotator beliefs followed by more recent models by Raykar et al. [7] and Zhou et al. [29]. Audhkhasi et al. [8] further extended the model to account for diversity in the reliability of annotators over the feature space. Our algorithm carries similar goals as Chen et al. [26] and Kumar et al. [27] to fuse preferences from multiple annotators, with modeling schemes inspired from proposals by Raykar et al. [7] and Audhkhasi et al. [8]. In the next section, we discuss the algorithms designed for the fusion of noisy pairwise comparisons from multiple annotators along with a few other baseline methods.

### 3 Methodology

Given a set of  $N$  items  $\mathbf{O} = \{O_1, O_2, \dots, O_N\}$  and  $K$  annotators, we represent the  $k^{\text{th}}$  annotator’s preference of  $O_i$  over  $O_j$  as  $O_i^k \succ O_j^k$ . Our goal is to infer the latent ground truth denoted by  $O_i \succ O_j$ , indicating that  $O_i$  is ranked higher than  $O_j$ . We also assume the availability of attributes/feature values  $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$  for each of the  $N$  objects, where  $\mathbf{x}_i$  is a vector of attributes for the item  $O_i$ . We define the variables  $z_{ij}^k$  ( $k = 1..K$ ) and  $z_{ij}^*$  to represent the preferences of the annotators and the ground truth as follows.

$$z_{ij}^* = \begin{cases} 1 & \text{if } O_i \succ O_j \\ 0 & \text{if } O_j \succ O_i \end{cases} \text{ and, } z_{ij}^k = \begin{cases} 1 & \text{if } O_i^k \succ O_j^k \\ 0 & \text{if } O_j^k \succ O_i^k \end{cases}, k = 1..K \quad (1)$$

Below we describe four methods to obtain the ground truth given the noisy pairwise comparisons between items. The first two methods, majority vote and Independent Annotator Model serve as a baseline. Although fusion from these methods is easy to perform, they assume that each annotator is equally reliable in inferring the ground truth which may not always be the case. The next two methods, the Joint Annotator Modeling and Variable

Reliability Joint Annotator Modeling schemes learns a reliability metric for each annotator. The final decision is made based on available annotations as well as the attributes for the pair of objects at hand.

### 3.1 Majority Vote (MV)

Majority voting is one of the most popular methods for merging decisions from multiple annotators and has been consistently used in various classification experiments [30, 31] as well as ranking [27]. In this method, we say that the inferred preference is  $O_i \succ O_j$  if a majority of the annotators say so. In case of a tie among annotators, a random decision is taken between  $z_{ij}^* = 1$  and  $z_{ij}^* = 0$ . Note that this model does not use the object attributes  $\mathbf{X}$  in inferring  $z_{ij}^*$  and relies solely on  $z_{ij}^k$  as shown in the graphical model in Figure 1(a). Also, each annotator is weighted equally in deciding the majority.

### 3.2 Independent Annotator Modeling (IAM)

In this scheme, we initially train annotator specific ranking models to capture the relation between object attributes and each annotator’s preference rankings. The ranking model for the  $k^{\text{th}}$  annotator returns a score  $f_k(\mathbf{x}_i)$  for every object  $O_i$  based on the attributes  $\mathbf{x}_i$ . Finally, the inferred ground truth value for  $z_{ij}^*$  is given by comparing the sum of scores  $f_k(\mathbf{x}_i)$  and  $f_k(\mathbf{x}_j)$  over all the annotators ( $k = 1..K$ ). This method is synonymous to the Bradley-Terry model [32] (extended by Chen et al. [26]) and the Borda count method [9] used for aggregating decisions from multiple annotators. In case of Bradley-Terry model, preference between two objects is determined based on their relevance scores, computed as a sum of  $f_k(\mathbf{x}_j)$  over all the annotators in the current IAM scheme. Similarly, in the Borda count method, each annotator scores every object and  $z_{ij}^*$  is inferred by comparing sum of scores across all the annotators. In this section, our substitute for the Borda count score for  $O_i$ , as given by the annotator  $k$ , is the value  $f_k(\mathbf{x}_i)$ . We describe the model training and ground truth inference in detail below.

**Training annotator specific models:** Given the  $k^{\text{th}}$  annotator’s pairwise preferences  $z_{ij}^k$ , we train an annotator specific Support Vector Ranker (SVR) [33] as the function  $f_k$ . Our goal is to learn  $f_k$  for every annotator  $k$ , such that the following holds.

$$O_i^k \succ O_j^k \iff z_{ij}^k = 1 \iff f_k(\mathbf{x}_i) > f_k(\mathbf{x}_j) \quad (2)$$

In this work, we chose  $f_k$  to be a linear function characterized by a weight vector  $\mathbf{w}_k$  such that  $f_k(\mathbf{x}_i) = \langle \mathbf{w}_k, \mathbf{x}_i \rangle$ , where  $\langle \mathbf{w}_k, \mathbf{x}_i \rangle$  represents the

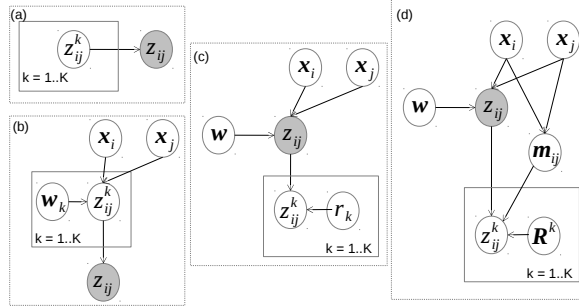


Figure 1: Graphical models for (a) Majority vote (MV) (b) Independent Annotator Model (IAM) (c) Joint Annotator Model (JAM) and, (d) Variable Reliability Joint Annotator Model (VRJAM) schemes.

dot product between  $w_k$  and  $x_i$ . An SVR targeting the problem in (2) performs the following optimization on the cost function  $\mathcal{M}_k$  [33].

$$w_k = \arg \min_{w_k} \mathcal{M}_k = \arg \min_{w_k} \sum_{\substack{\text{All pairs} \\ x_i, x_j}} z_{ij}^k [1 - \langle w_k, \{x_i - x_j\} \rangle]_+ + (1 - z_{ij}^k) [1 - \langle w_k, \{x_j - x_i\} \rangle]_+ \quad (3)$$

In the equation above,  $\{x_i - x_j\}$  depicts a notion of difference operator between  $x_i$  and  $x_j$  and  $[\ ]_+$  represents the standard hinge loss function [34]. In this work, we use  $\{x_i - x_j\}$  to be a simple element-wise subtraction between attribute vectors  $x_i$  and  $x_j$ . We learn  $w_k$  ( $\forall k = 1..K$ ) using the standard gradient descent algorithm [35]. Since  $\mathcal{M}_k$  is non differentiable, we use the approximation suggested by Rennie et al. [36] in the hinge loss function.

**Fusing annotator models:** After obtaining  $f_k$  for each of the annotators, we say  $z_{ij} = 1$  if:

$$\sum_{k=1}^K f_k(x_i) > \sum_{k=1}^K f_k(x_j) \quad (4)$$

A graphical model representing this scheme is shown in Figure 1(b). Note in order to obtain  $z_{ij}^*$ , the scheme of unweighted combination is enforced on  $f_k$  outputs.

### 3.3 Joint Annotator Modeling (JAM)

In this section, we propose an Expectation-Maximization (EM) algorithm [11] to infer the ground truth by jointly modeling the noisy comparisons. Our algorithm is inspired by similar works [7, 8] in the domain of

classification problems. A graphical model for this scheme is shown in Figure 1(c). We assume the ground truth  $z_{ij}^*$  to be a latent variable that can be inferred using the object attributes  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . Our choice for inferring  $z_{ij}^*$  based on  $\mathbf{x}_i, \mathbf{x}_j$  is again an SVR model with a weight vector  $\mathbf{w}$ . Furthermore, we assume that  $z_{ij}^k$  is obtained by flipping the binary variable  $z_{ij}^*$  with a probability  $r_k$ . In summary, this model assumes that there is an inherent true preference given attributes from two objects and the annotators are flipping it based on annotator specific probabilities ( $r_k, k = 1..K$ ). Consequently, the probability  $r_k$  also provides a measure of annotator quality as a higher  $r_k$  implies higher chances of an annotator committing an error. We infer the latent ground truth  $z_{ij}^*$  using an EM algorithm described in the next section.

### 3.3.1 Expectation-Maximization algorithm

The EM algorithm maximizes the log-likelihood  $\mathcal{L}$  of the observed data, that is, annotator preferences given the object attributes and the model parameters. In our case,  $\mathcal{L}$  is given as shown in (5). Notice the introduction of the latent ground truth  $z_{ij}^*$  into  $\mathcal{L}$  in (6).

$$\mathcal{L} = \sum_{\substack{\text{All pairs} \\ \mathbf{x}_i, \mathbf{x}_j}} \log p(z_{ij}^1, \dots, z_{ij}^K / \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) \quad (5)$$

$$= \sum_{\substack{\text{All pairs} \\ \mathbf{x}_i, \mathbf{x}_j}} \sum_{z_{ij}^*} \log p(z_{ij}^*, z_{ij}^1, \dots, z_{ij}^K / \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) \quad (6)$$

Following the EM derivation procedure in section 9.4 in [12], we introduce a distribution over the latent ground truth  $z_{ij}^*$ :  $q(z_{ij}^*)$ . Consequently  $\mathcal{L}$  can be written as sum of two terms, a Kullback Leibler (KL) divergence term  $\text{KL}(q||p)$  and another log-likelihood term  $\mathcal{M}$  as shown in (7).

$$\mathcal{L} = \mathcal{M} + \text{KL}(q||p) \quad (7)$$

where,

$$\begin{aligned} \mathcal{M} = & \sum_{\text{all pairs } \mathbf{x}_i, \mathbf{x}_j} \sum_{z_{ij}^*} q(z_{ij}^*) \times \\ & \log \left\{ \frac{p(z_{ij}^*, z_{ij}^1, \dots, z_{ij}^K | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K)}{q(z_{ij}^*)} \right\} \end{aligned} \quad (8)$$

$$\text{KL}(q||p) = - \sum_{\text{all pairs } \mathbf{x}_i, \mathbf{x}_j} \sum_{z_{ij}^*} q(z_{ij}^*) \times \log \left\{ \frac{p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K)}{q(z_{ij}^*)} \right\} \quad (9)$$

The EM algorithm consists of two steps: the E and M steps. In the E-step,  $\mathcal{M}$  is maximized with respect to  $q(z_{ij}^*)$  while holding the other parameters constant. The solution is equivalent to the posterior distribution  $p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K)$ . In the M-step,  $\mathcal{M}$  is maximized with respect to model parameters while holding the estimated distribution  $q(z_{ij}^*)$  constant. We describe the parameter initialization followed by the E and M steps below.

**Initialization:** We randomly initialize the SVR weight vector  $\mathbf{w}$  and the probabilities of flipping  $r_k (k = 1..K)$ .

**While**  $\mathbf{w}, r_1, \dots, r_K$  not converged perform E and M-steps, where:

E-step: In the E-step, we set the probability distribution  $q(z_{ij}^*)$  equal to  $p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K)$ . This quantity can be represented as shown in (10). A detailed derivation for this quantity can be seen in Appendix 1.

$$q(z_{ij}^*) = p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) = \left( p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \times \prod_{k=1}^K p(z_{ij}^k | z_{ij}^*, r_k) \right) / p(z_{ij}^1, \dots, z_{ij}^K) \quad (10)$$

Note that the first term  $p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w})$  in (10) is conditioned on the SVR model parameters  $\mathbf{w}$  and object attributes  $\mathbf{x}_i$  and  $\mathbf{x}_j$  only. Since SVR is not a probabilistic model, we apply a commonly used trick in support vector machine classifiers employed to obtain class probabilities. The trick involves fitting logistic models to distance from the decision hyperplane to obtain the probabilities of preference decisions [37] (A comparison of the hinge loss function and the logistic loss function is made in Appendix 3). Equations (11) and (12) show the computation for  $p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w})$  using the logistic model.

$$p(z_{ij}^* = 1 | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) = \frac{\exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle}{1 + \exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle} \quad (11)$$

$$p(z_{ij}^* = 0 | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) = 1 - p(z_{ij}^* = 1 | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \quad (12)$$



The second term  $p(z_{ij}^k|z_{ij}^*, r_k)$  in (10) is  $r_k$  if  $z_{ij}^k$  and  $z_{ij}^*$  are in disagreement and  $1 - r_k$  otherwise, as shown below.

$$p(z_{ij}^k|z_{ij}^*, r_k) = \begin{cases} r_k & \text{if } z_{ij}^k \neq z_{ij}^* \\ 1 - r_k & \text{if } z_{ij}^k = z_{ij}^* \end{cases} \quad (13)$$

Replacing the values in (10) from (11) and (13), we can represent  $q(z_{ij}^* = 1)$  as shown in (14).  $q(z_{ij}^* = 0)$  can be computed accordingly.

$$q(z_{ij}^* = 1) = \frac{\exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle}{1 + \exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle} \times \prod_{k=1}^K \underbrace{[(r_k)^{(1-z_{ij}^k)} \times (1-r_k)^{z_{ij}^k}]}_{r_k/(1-r_k) \text{ is multiplied when } z_{ij}^k = 0/(1)} / p(z_{ij}^1, \dots, z_{ij}^K) \quad (14)$$

Note that the denominator  $p(z_{ij}^1, \dots, z_{ij}^K)$  is common between  $q(z_{ij}^* = 1)$  and  $q(z_{ij}^* = 0)$  and need not be computed. We can just compute the numerator in (10) for  $q(z_{ij}^* = 1)$  and  $q(z_{ij}^* = 0)$  and normalize these probabilities to sum to one. In the next section, we discuss the M-step.

M-step: In this step, we estimate the model parameters  $\mathbf{w}, r_k (k = 1..K)$  based on estimated distribution  $q(z_{ij}^*)$ . These parameters are estimated by maximizing  $\mathcal{M}$  after substituting  $q(z_{ij}^*)$  estimated in the E-step. In our case,  $\mathcal{M}$  can be written as shown in (15).  $\mathbb{H}(q(z_{ij}^*))$  is the entropy of  $q(z_{ij}^*)$  and is a constant term with respect to the model parameters  $\mathbf{w}, r_1, \dots, r_K$ . We disregard the entropy term for further M-step derivations.

$$\mathcal{M} = \sum_{z_{ij}^*} q(z_{ij}^*) \times \log p(z_{ij}^*, z_{ij}^1, \dots, z_{ij}^K | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) + \mathbb{H}(q(z_{ij}^*)) \quad (15)$$

We can rewrite  $\mathcal{M}$  as shown in (16). For a detailed derivation, please see Appendix 2. Note that each parameter  $\mathbf{w}, r_1, \dots, r_K$  appears in a separate term within the summation in (16) and thus, we only need to consider the corresponding term while optimizing for a parameter. We discuss the optimization for the SVR parameters  $\mathbf{w}$  and flipping probabilities  $r_k$  below.

$$\mathcal{M} = \sum_{z_{ij}^*} q(z_{ij}^*) \left( \sum_{k=1}^K \log p(z_{ij}^k | z_{ij}^*, r_k) + \log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \right) \quad (16)$$

*Obtaining SVR weight vector  $\mathbf{w}$ :* We only need to consider the following term  $\mathcal{M}_{\mathbf{w}}$  within  $\mathcal{M}$  to optimize for  $\mathbf{w}$ .

$$\mathcal{M}_{\mathbf{w}} = \sum_{z_{ij}^*} q(z_{ij}^*) \log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \quad (17)$$

In the EM algorithm,  $\log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w})$ , would be obtained from a probabilistic model to infer  $z_{ij}^*$  conditioned on  $\mathbf{x}_i, \mathbf{x}_j, \mathbf{w}$ . However, since the choice of our model is a non-probabilistic SVR, we instead solve the following optimization in (18) to obtain  $\mathbf{w}$ . We would like to point out that this is an approximation we use in the EM algorithm. Appendix 3 shows the probability distribution corresponding to the logistic models used in (11) and its relation to the following optimization.

$$\begin{aligned} \mathbf{w} = \arg \min_{\mathbf{w}} \mathcal{M}'_{\mathbf{w}} = \arg \min_{\mathbf{w}} & \left( q(z_{ij}^* = 1) [1 - \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle]_+ \right. \\ & \left. + q(z_{ij}^* = 0) [1 - \langle \mathbf{w}, \{\mathbf{x}_j - \mathbf{x}_i\} \rangle]_+ \right) \end{aligned} \quad (18)$$

Note that  $\mathcal{M}'_{\mathbf{w}}$  in (18) is similar to the cost function  $\mathcal{L}_k$  defined in (3) for training annotator specific models. However, instead of being trained on binary decisions values (e.g.,  $z_{ij}^k$  used in  $\mathcal{L}_k$ ),  $\mathcal{M}_{\mathbf{w}}$  is defined over the soft estimate  $q(z_{ij}^*)$ . Next, we discuss the optimization problem to obtain  $r_1, \dots, r_K$ .

*Obtaining probability of flipping  $r_1, \dots, r_K$ :* In order to obtain  $r_k$ , we need to optimize the following term within  $\mathcal{M}$ .

$$\begin{aligned} r_k &= \arg \min_{r_k} \mathcal{M}_r^k \\ &= \arg \min_{r_k} \sum_{\text{All pairs } \mathbf{x}_i, \mathbf{x}_j} \sum_{z_{ij}^*} q(z_{ij}^*) \log p(z_{ij}^k | z_{ij}^*, r_k) \end{aligned} \quad (19)$$

$p(z_{ij}^k | z_{ij}^*, r_k)$  is replaced in the above equation as shown in (13) and the term can be optimized to obtain  $r_k$ . We obtain the final inference for  $z_{ij}^*$  as discussed below.

**Final inference:** After convergence, we make the final inference on  $z_{ij}^*$  based on obtained distribution  $p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K)$ , as was derived in (10)-(14).  $z_{ij}^*$  is inferred to be 1 or 0 as per the following equation.

$$\begin{aligned} p(z_{ij}^* = 1 | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) &\stackrel{1}{>}_0 \\ p(z_{ij}^* = 0 | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) &\end{aligned} \quad (20)$$

Next, we propose a modification to this scheme considering the annotators' probability of flipping to be variable.

### 3.4 Variable Reliability Joint Annotator Modeling (VRJAM)

This scheme is similar to the joint annotator model proposed in the previous section except for the probability of flipping  $r_k$  being variable. The motivation behind this scheme is that annotators may have variable reliability depending upon the pair of objects  $O_i$  and  $O_j$  at hand (a similar assumption is in the model proposed by Audhkhasi et al. [8]). Therefore, instead of a constant  $r_k$  for the annotator  $k$ , we determine a vector  $\mathbf{R}_k = [r_k^1, \dots, r_k^D]$ , where based on the difference vector  $\{\mathbf{x}_i - \mathbf{x}_j\}$ , one of the values  $r_k^d (d = 1..D)$  is chosen as the probability of flipping. We retain the assumption that  $z_{ij}^*$  is a latent variable conditioned on  $\mathbf{x}_i, \mathbf{x}_j$  and the SVR weight vector  $\mathbf{w}$ . We again train this model using an EM algorithm described below. The algorithm is similar to the EM algorithm in section 3.3 and we borrow several steps for the sake of brevity.

#### 3.4.1 Expectation-Maximization algorithm

For the purpose of our experiments, we divide the space spanned by difference vectors  $\{\mathbf{x}_i - \mathbf{x}_j\}$  into  $D$  clusters. For the  $k^{\text{th}}$  annotator, a distinct probability of flipping  $r_k^d (d = 1..D)$  is computed in each cluster. We obtain the clusters by performing the standard K-Means clustering [38] on the values  $\{\mathbf{x}_i - \mathbf{x}_j\}$  obtained over all pairs  $\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}$ . The membership of  $\{\mathbf{x}_i - \mathbf{x}_j\}$  to a cluster is denoted by a 1-in- $D$  encoding vector  $\mathbf{m}_{ij} = [m_{ij}^1, \dots, m_{ij}^D]$  where  $m_{ij}^d = 1$  indicates that  $\{\mathbf{x}_i - \mathbf{x}_j\}$  belongs to the  $d^{\text{th}}$  cluster. The overall graphical model for this scheme is represented in Figure 1(d). The graphical model is very similar to the one in Figure 1(c), except for  $\mathbf{m}_{ij}$  now determining the flipping probability. The data log-likelihood term  $\mathcal{L}$  in (5) changes slightly to incorporate  $\mathbf{R}_1, \dots, \mathbf{R}_K$  and  $\mathbf{m}_{ij}$  (instead of scalar values  $r_1, \dots, r_K$ ) as represented by  $\mathcal{L}'$  in (21). We perform the initialization, the E and M-steps and final inference as discussed in the next section.

$$\mathcal{L}' = \sum_{\substack{\text{All pairs} \\ \mathbf{x}_i, \mathbf{x}_j}} \log p(z_{ij}^1, \dots, z_{ij}^K / \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, \mathbf{R}_1, \dots, \mathbf{R}_K, \mathbf{m}_{ij}) \quad (21)$$

**Initialization:** We randomly initialize the SVR weight vector  $\mathbf{w}$  and the vectors  $\mathbf{R}_k$  for all the annotators. We perform K-means clustering to segment the space spanned by  $\{\mathbf{x}_i - \mathbf{x}_j\}, \forall \mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}$ . The number of clusters  $D$  is set empirically by gradually increasing  $D$  until the distance between two cluster centroids falls below a threshold (compared to distances to other centroids).

While  $\mathbf{w}, \mathbf{R}_1, \dots, \mathbf{R}_K$  not converged perform E and M-steps, where:

E-step: The E-step is same as the E-step in section 3.3. The only difference is that  $p(z_{ij}^k | z_{ij}^*, r_k)$  in (10) is replaced by  $p(z_{ij}^k | z_{ij}^*, \mathbf{R}_k, \mathbf{m}_{ij})$ , which equals to the quantity in (22).  $\langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle$  represents a dot product between  $\mathbf{R}_k$  and  $\mathbf{m}_{ij}$  to select an entry in  $\mathbf{R}_k$  based on the cluster index corresponding to  $\{\mathbf{x}_j - \mathbf{x}_i\}$ .

$$p(z_{ij}^k | z_{ij}^*, \mathbf{R}_k, \mathbf{m}_{ij}) = \begin{cases} \langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle, & \text{if } z_{ij}^k \neq z_{ij}^* \\ 1 - \langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle, & \text{if } z_{ij}^k = z_{ij}^* \end{cases} \quad (22)$$

Consequently,  $q(z_{ij}^* = 1)$  is computed as shown in (23). After estimating  $q(z_{ij}^*)$ , we estimate the model parameters as discussed next.

$$q(z_{ij}^* = 1) = \frac{\exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle}{1 + \exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle} \times \prod_{k=1}^K \underbrace{[\langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle^{(1-z_{ij}^k)} \times (1 - \langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle)^{z_{ij}^k}]}_{\langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle / (1 - \langle \mathbf{R}_k, \mathbf{m}_{ij} \rangle) \text{ is multiplied when } z_{ij}^k = 0(1)} / p(z_{ij}^1, \dots, z_{ij}^K) \quad (23)$$

M-step: In the M-step, we re-estimate the parameters  $\mathbf{w}$  and the vectors  $\mathbf{R}_k$ . Value of  $\mathcal{M}$  also alters in this formulation to incorporate  $\mathbf{R}_1, \dots, \mathbf{R}_K$  and  $\mathbf{m}_{ij}$ .  $p(z_{ij}^k | z_{ij}^*, r_k)$  in (16) is replaced by  $p(z_{ij}^k | z_{ij}^*, \mathbf{R}_k, \mathbf{m}_{ij})$ . This has no impact on the estimation of  $\mathbf{w}$ , which remains the same as in section 3.3. We describe the estimation of the vector  $\mathbf{R}_k$  below.

*Obtaining probability of flipping entries in  $\mathbf{R}_k$* : The optimization framework to obtain  $\mathbf{R}_k$  is shown below.

$$\mathbf{R}_k = \arg \min_{\mathbf{R}_k} \sum_{\substack{\text{All pairs} \\ \mathbf{x}_i, \mathbf{x}_j}} \sum_{z_{ij}^*} q(z_{ij}^*) \log p(z_{ij}^k | z_{ij}^*, \mathbf{R}_k, \mathbf{m}_{ij}) \quad (24)$$

The above optimization over the vector  $\mathbf{R}_k$  can easily be broken down into scalar optimization over each of its entries after replacing  $p(z_{ij}^k | z_{ij}^*, \mathbf{R}_k, \mathbf{m}_{ij})$  as shown in (22). We next discuss the final step for inferring  $z_{ij}^*$ .

**Final inference**: The final inference on  $z_{ij}^*$  is made based the following likelihood comparison once the model converges. This inference is similar to one in the JAM scheme.

$$\begin{aligned}
p(z_{ij}^* = 1 | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, \mathbf{R}_1, \dots, \mathbf{R}_K, m_{ij}) & \stackrel{1}{<} \\
p(z_{ij}^* = 0 | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, \mathbf{R}_1, \dots, \mathbf{R}_K, m_{ij}) & \stackrel{0}{>}
\end{aligned} \tag{25}$$

In the next section, we evaluate various fusion schemes on several datasets with synthetic annotations as well as annotations obtained from machines and humans.

## 4 Experimental Results

We test the discussed ranking algorithms on two synthetically created data sets and two real world data set as discussed next.

### 4.1 Data sets with synthetic annotations

We use the two wine quality data sets (red and white wine data sets) [39] available in the UCI data repository [40]. Each data set provides 11 attributes for each entry and a quality score between 0-10 (10 being the best). In pairwise comparison between two entries  $O_i$  and  $O_j$ , we say that the ground truth is  $z_{ij} = 1$  if  $O_i$  has a higher quality score than  $O_j$ . Below we provide a short description of synthetic creation of noisy annotator labels from this data set followed by a set of three experiments investigating the reliability inference for each annotator and the effect of quality and number of annotators.

**Creating synthetic noisy annotations:** Given the number of annotators  $K$ , we create synthetic noisy annotations for the  $k^{\text{th}}$  annotator by flipping the ground truth  $z_{ij}^*$  based on a Bernoulli variable. The parameter of the Bernoulli variable for annotator  $k$  is denoted by  $b_k$  and a higher  $b_k$  implies higher chances of  $z_{ij}^*$  being flipped. In the first experiment presented in the next section, we investigate the relation between  $b_k$  used for each annotator and the probability of flipping  $r_k$  determined by our joint annotator models.

#### 4.1.1 Relationship between probability of flipping and annotator noise

In this experiment, we use a set of 6 noisy annotators with  $b_k = k/20$ . That is the first annotator is the best annotator with only 5% chance of flipping where as the sixth annotator has a 30% chance of flipping. We train the Joint Annotator Model (JAM) and Variable Reliability Joint Annotator Model (VRJAM). Table 1 shows the values for  $r_k$  estimated using JAM and the mean value of vector  $\mathbf{R}_k$  estimated using VRJAM on the red wine data set (similar patterns are observed for white wine data set). Higher values for  $r_k$  and mean of  $\mathbf{R}_k$  imply that the annotator  $k$  is inferred to be more

Table 1: Values of  $r_k$  &  $\text{mean}(\mathbf{R}_k)$  obtained on the red wine data set.

Model	Parameter	Values for $k = 1..6; b_k = k/20$
JAM	$r_k$	$\{.032, .086, .176, .196, .246, .273\}$
VRJAM	$\text{Mean}(\mathbf{R}_k)$	$\{.033, .085, .175, .196, .245, .273\}$

Table 2: Accuracy in inferring  $z_{ij}^*$  in the synthetic data sets.

Data set	MV	IAM	JAM	VRJAM
Red wine	95.9	55.2	97.9	98.0
White wine	96.1	55.3	97.9	98.1

noisy. We also show the model accuracy in inferring the ground truth  $z_{ij}^*$  over all pairs of objects in the data set in Table 2.

From the Table 1, we observe that as the noise increases over annotators, our model successfully infers a higher probability of flipping. The values  $r_k$  and the mean of vector  $\mathbf{R}_k$  are fairly close to each other indicating that the JAM and VRJAM model are very similar in inferring probability of flipping. This is expected as VRJAM differs from JAM only in determining cluster-wise probabilities and their average should be fairly close to  $r_k$ . From Table 2, we observe that the proposed models outperform Majority Vote (MV) and Independent Annotator Modeling (IAM). The difference in performance of JAM and VRJAM is not significant. This stems from the choice of synthetic annotation generation as the noise added to the annotations is uniform and does not change based on the pair of objects at hand. Therefore VRJAM has no particular modeling advantage over the JAM scheme. Also, the performance of IAM is particularly low. Our investigation reveals that the performances of the individual annotator SVR models ( $f_k$  in section 3.2) were very low (e.g., varied between 53.0%-64.4% in red wine data set). Since IAM performs a sum of  $f_k$  over these fairly weak models, the final performance is poor. This shows that the IAM performance is contingent upon the model choice and can improve with a better choice for  $f_k$ . However, an interesting point to note here is that the IAM performance (e.g., 55.2% for red wine data set) lies between the performance of the best annotator (64.4% for red wine data set) and the worst annotator (53.0% for red wine data set). This reflects the fact that IAM is susceptible to performing below collective knowledge of the crowd and can perform worse than the best available annotator.

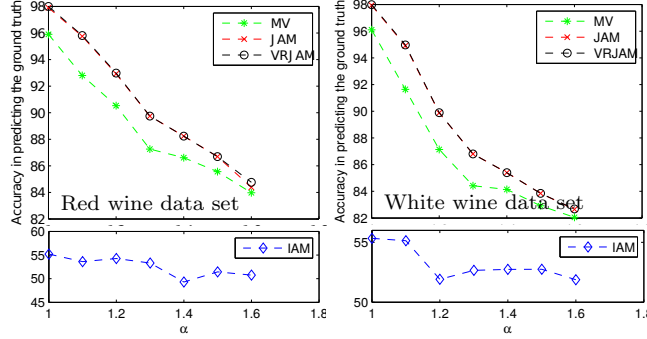


Figure 2: Model performances with increasing annotator noises.

#### 4.1.2 Relationship between model performances and annotator noise

In this section, we perform multiple experiments similar to the one mentioned in the previous section. We chose a set of 6 annotators and in each experiment, we increase the parameter  $b_k$ . Within an experiment,  $b_k$  for the annotator  $k$  is set at  $\alpha k/20$  and the parameter  $\alpha$  is increased by 10% over consecutive experiments. We plot the accuracy of the MV, IAM, JAM and VRJAM algorithms in inferring the ground truth  $z_{ij}^*$  with increasing  $\alpha$  in Figure 2. From the figure, we note that the model performance drops as the annotator noise increases. Performances of the VRJAM and JAM schemes are again similar because of the reasons stated in the previous section. Another interesting observation is that the performances of MV, JAM and VRJAM converge as the annotator noise increases. This indicates that the joint models are likely to perform better than MV with better quality annotators. The IAM performances are again low attributed to weak annotator modeling by the SVRs.

#### 4.1.3 Relationship between model performances and number of annotators

In this section, we perform multiple experiments by varying the total count of annotators  $K$ . The parameter  $b_k$  for the annotator  $k$  is kept constant at  $k/20$ . Figure 3 shows the plots for model performance as  $K$  is varied from 3 to 9. In this case, we observe that except for IAM, performance of all models increase with increase in number of annotators. This indicates that addition of more noisy annotators (as  $b_k < b_{k+1}$ ) tends to decrease IAM performance. Also, the JAM and VRJAM models provide greater improvement over MV with addition of more annotators. The performance of MV, JAM and VRJAM models are same at  $K = 3$  and the absolute improvement

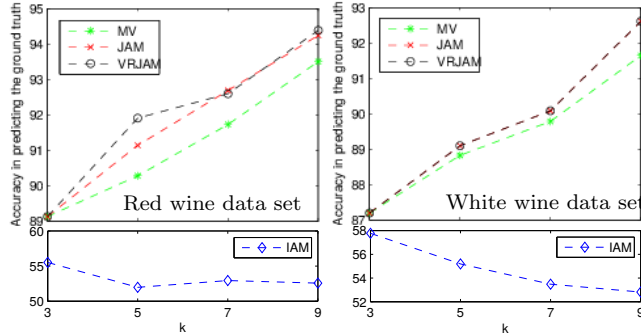


Figure 3: Model performances with increasing number of annotators.

of the joint models over MV increases as we add more annotators. VRJAM and JAM again perform at similar levels. As stated, we attribute this to the nature of our synthetic labels creation where noisy annotators flip  $z_{ij}^*$  solely based on  $b_k$  and not based on the object attributes  $\mathbf{x}_i, \mathbf{x}_j$ . In the next section, we test our algorithms on a data sets with machine/human annotations and analyze the results.

## 4.2 Data set with machine/human annotations

We show the results for two real world data sets, one annotated by machine experts and other by naive mechanical turk workers. We discuss the results for these two datasets below.

### 4.2.1 Digit ranking dataset: Machine annotation

We use a subset of the pen based recognition of handwritten digits dataset [41] to rank images based on the digit value contained (for instance image with digit 9 is ranked higher than image containing any other digit). The dataset contains 1k samples of images with 16 features, leading to 370k possible comparisons (we do not consider comparison between images containing same values). We initially annotate the dataset using a set of five classifier as machine annotators: K-Nearest Neighbors (KNN), Logistic Regression, Naive Bayes, Random Forest and Perceptron [12]. These annotations are obtained using a 10 fold cross-validation framework. Each classifier is trained on a subset of 3-4 features (out of 16) on 90% of the data and results are obtained on the remaining 10%. This process is repeated till we annotate the entire data using the classifiers. Note that in this dataset, we have access to the ground truth which may not always be the true (this is the case with the dataset in the next section). Table 3 shows the performance of each classifier as a machine annotator in pairwise comparison between images. Table 4 shows the performance of the fusion schemes operating over



Table 3: Ratio of pairwise comparisons in which a classifier ranks the image containing greater value higher than the other image in the pair. (KNN: KNN classifier, LR: Logistic Regression, NB: Naive Bayes classifier, RF: Random Forests classifier and Perc.: Perceptron).

Classifier	KNN	LR	NB	RF	Perc.
Performance	67.8	69.1	69.0	72.0	59.9

Table 4: Performance of the fusion schemes on pairwise comparisons  $z_{ij}^k$ , as obtained from the machine annotators.

Fusion scheme	MV	IAM	JAM	VRJAM
Performance	78.0	65.9	78.1	79.7

the machine annotations thus obtained. We use the entire set of 16 features in the JAM and VRJAM fusion schemes.

From the Table 3, we see that the machine annotators perform in the range of 59% to 72% on the metric of pairwise comparison accuracy. Results in Table 4 indicate that the MV, JAM and VRJAM schemes outperform the best machine annotator, i.e., random forests. Where as the performances of MV and JAM are not significantly different, VRJAM performs significantly better than both MV and JAM schemes (McNemar’s test [42], significance level: 5%, computed over the 370k comparison samples). This indicates that assigning a flipping probability conditioned on the pair of images at hand is essential in this data set. The IAM scheme again fails to beat the best annotator and performs at a value within the range of best and the worst annotator. This indicates that an unweighted fusion of experts may perform below the collective knowledge of the crowd and weighting annotators based on individual performances may help. In the next section, we test the fusion scheme on another real data set with human annotators.

#### 4.2.2 Safari Bob dataset

In this section, we test our algorithms on the Safari Bob data set [43]. This data set involves two populations of High Functioning Autism (HFA) and Typically Developing (TD) individuals retelling a story based on a video stimulus. The recording of story retelling are later rated by naive Mechanical Turk (MTurk) raters for expressiveness and naturalness on a scale from 0-4 (4 being the best). We use a set of 40 TD kids and 65 HFA kids rated by 5 annotators and infer the ground truth expressiveness and naturalness from the available ratings. The attributes  $\mathbf{x}_i$  we use to train the models are statistical functionals extracted on prosodic and spectral features from the kid’s speech (mean and variance of pitch, intensity, Mel filter banks and Cepstral

Table 5: Comparison of expressiveness/naturalness between TD and HFA kids. TD kids are expected to be more expressive/natural.

Attribute	Ratio of times TD kids are inferred to have a higher rank over HFA kids			
	MV	IAM	JAM	VRJAM
Expressiveness	64.3	61.5	64.3	65.4
Naturalness	55.7	52.7	55.9	57.7

Coefficients) as are also used in [30, 43]. Since we do not have the ground truth available for evaluation, we analyze the association of inferred expressiveness and naturalness with the population attributes of HFA and TD. Although the relationship between autism and expressiveness/naturalness is fairly complex and undergoing extensive investigation [44], TD kids are expected to be ranked higher in expressiveness/naturalness over HFA kids [43]. We infer the latent ground truth for expressiveness/naturalness using our models set and show (Table 5) the proportion of times the models infer TD kids to have a higher expressiveness/naturalness than HFA kids.

From the results, we observe that a TD kid is more often inferred to have a higher expressiveness/naturalness over a HFA kid. Whereas outputs for MV and JAM are fairly close to each other, the outputs from the VRJAM has the highest proportion of times that a TD kid is inferred to be more expressive/natural than an ASD kid. This trend is encouraging although the relation between speech expressiveness/naturalness and autism may not be this straightforward. Due to unavailability of ground truth, this experiment can not be used to support the efficacy of proposed algorithms. However the observed results motivate the application of proposed algorithms to data sets where the ground truth is unobserved.

Overall, the experiment on synthetic, machine and human annotations in this section provide an understanding of the proposed algorithms within the aspects of annotator reliability, quality, and number of annotators. Although the performance of VRJAM is not significantly better in the case of synthetic annotations, results on the machine and human annotations indicate the importance of accounting for differences in the reliability of annotators based on the pair of objects at hand. We conclude our work in the next section and present a few future directions.

## 5 Conclusion

In this paper, we address the problem of inferring the hidden ground truth preference given noisy annotations from multiple annotators. We propose an EM algorithm based Joint Annotator Modeling (JAM) scheme, con-

sidering the latent ground truth preference to be a hidden variable and inferring it based on available annotation and object attributes. Given a pair of objects, the JAM scheme infers the latent true preference order based on a set of object attributes as well as noisy annotator preferences. The model assumes that annotators flip the true preference order based on a Bernoulli random variable and estimates annotator specific “probability of flipping”. We further extend the model to estimate a non-constant “probability of flipping” conditioned on the pair of objects at hand in the Variable Reliability Joint Annotator Model (VRJAM). We test the JAM and VRJAM schemes against majority voting and Independent Annotator Modeling schemes on data sets with annotations obtained synthetically, from machines as well as from human annotators. Using the data set with synthetic annotations, we test the impact of annotator quality and quantity on our models. The results on data sets with machine annotations depicts the importance of having a variable reliability per annotator based on pair of objects at hand. Finally, in the Safari Bob data with human annotators, we interpret the results based on the expected trends of expressiveness/naturalness in TD and HFA kids.

In the future, we aim to extend the presented algorithms by integrating other existing work in the ranking domain (e.g., active learning). Other work in rank aggregation inferring a rank order probability distribution can also be integrated into the proposed EM framework. Also, within classification there are further extensions of multiple annotator models which can be incorporated into the current EM framework. We also aim to implement the designed algorithms to other data sets such as the Safari Bob data set in understanding the diversity in perception of various psychological constructs (e.g. naturalness) by the human annotators and their relation to a target variable (e.g. autism severity).

## References

- [1] E. Hüllermeier, J. Fürnkranz, W. Cheng, and K. Brinker, “Label ranking by learning pairwise preferences,” *Artificial Intelligence*, vol. 172, no. 16, 2008.
- [2] T.-Y. Liu, “Learning to rank for information retrieval,” *Foundations and Trends in Information Retrieval*, vol. 3, 2009.
- [3] T. H. Haveliwala, “Topic-sensitive pagerank,” in *Proceedings of the 11th international conference on World Wide Web*. ACM, 2002.
- [4] K. H.-Y. Lin and H.-H. Chen, “Ranking reader emotions using pairwise loss minimization and emotional distribution regression,” in *Proceedings*

of the conference on empirical methods in natural language processing. Association for Computational Linguistics, 2008.

- [5] O. Wu, W. Hu, and J. Gao, “Learning to rank under multiple annotators,” in *IJCAI Proceedings-International Joint Conference on Artificial Intelligence*, 2011, vol. 22.
- [6] Y. Yan, R. Rosales, G. Fung, M. W. Schmidt, G. H. Valadez, L. Bogoni, L. Moy, and J. G. Dy, “Modeling annotator expertise: Learning when everybody knows a bit of something,” in *International conference on artificial intelligence and statistics*, 2010.
- [7] V. C. Raykar, S. Yu, L. H. Zhao, G. H. Valadez, C. Florin, L. Bogoni, and L. Moy, “Learning from crowds,” *The Journal of Machine Learning Research*, vol. 11, 2010.
- [8] K. Audhkhasi and S. Narayanan, “A globally-variant locally-constant model for fusion of labels from multiple diverse experts without using reference labels,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 4, 2013.
- [9] M. Van Erp and L. Schomaker, “Variants of the borda count method for combining ranked classifier hypotheses,” in *in the seventh international workshop on frontiers in handwriting recognition*. Citeseer, 2000.
- [10] E. M. Niou, “A note on nanson’s rule,” *Public Choice*, vol. 54, 1987.
- [11] T. K. Moon, “The expectation-maximization algorithm,” *Signal processing magazine, IEEE*, vol. 13, no. 6, 1996.
- [12] C. M. Bishop, *Pattern recognition and machine learning*, springer, 2006.
- [13] J. Fürnkranz and E. Hüllermeier, “Pairwise preference learning and ranking,” in *Machine Learning: ECML*. Springer, 2003.
- [14] K. Brinker, “Active learning of label ranking functions,” in *Proceedings of the twenty-first international conference on Machine learning*. ACM, 2004.
- [15] B. Long, O. Chapelle, Y. Zhang, Y. Chang, Z. Zheng, and B. Tseng, “Active learning for ranking through expected loss optimization,” in *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*. ACM, 2010.

- [16] W. Chu and Z. Ghahramani, “Extensions of gaussian processes for ranking: semisupervised and active learning,” in *Proceedings of the NIPS Workshop on Learning to Rank*. MIT, 2005.
- [17] J. He, M. Li, H.-J. Zhang, H. Tong, and C. Zhang, “Manifold-ranking based image retrieval,” in *Proceedings of the 12th annual ACM international conference on Multimedia*, 2004.
- [18] C. Quoc and V. Le, “Learning to rank with nonsmooth cost functions,” *Proceedings of the Advances in Neural Information Processing Systems*, vol. 19, 2007.
- [19] P. Li, Q. Wu, and C. J. Burges, “Mcrank: Learning to rank using multiple classification and gradient boosting,” in *Advances in neural information processing systems*, 2007.
- [20] K. Duh and K. Kirchhoff, “Learning to rank with partially-labeled data,” in *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*. ACM, 2008.
- [21] Z. Cao, T. Qin, T.-Y. Liu, M.-F. Tsai, and H. Li, “Learning to rank: from pairwise approach to listwise approach,” in *Proceedings of the 24th international conference on Machine learning*. ACM, 2007.
- [22] C. Dwork, R. Kumar, M. Naor, and D. Sivakumar, “Rank aggregation methods for the web,” in *Proceedings of the 10th international conference on World Wide Web*. ACM, 2001.
- [23] W. Ding, P. Ishwar, and V. Saligrama, “Learning shared rankings from mixtures of noisy pairwise comparisons,” in *Acoustics, Speech and Signal Processing (ICASSP), 2015 IEEE International Conference on*. IEEE, 2015.
- [24] D. F. Gleich and L.-h. Lim, “Rank aggregation via nuclear norm minimization,” in *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2011.
- [25] Y. Pan, H. Lai, C. Liu, Y. Tang, and S. Yan, “Rank aggregation via low-rank and structured-sparse decomposition,” in *AAAI*, 2013.
- [26] X. Chen, P. N. Bennett, K. Collins-Thompson, and E. Horvitz, “Pairwise ranking aggregation in a crowdsourced setting,” in *Proceedings of the sixth ACM international conference on Web search and data mining*. ACM, 2013.

- [27] A. Kumar and M. Lease, “Learning to rank from a noisy crowd,” in *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. ACM, 2011.
- [28] A. P. Dawid and A. M. Skene, “Maximum likelihood estimation of observer error-rates using the em algorithm,” *Applied statistics*, pp. 20–28, 1979.
- [29] D. Zhou, S. Basu, Y. Mao, and J. C. Platt, “Learning from the wisdom of crowds by minimax entropy,” in *Advances in Neural Information Processing Systems*, 2012, pp. 2195–2203.
- [30] R. Gupta, C.-C. Lee, and S. Narayanan, “Classification of emotional content of sighs in dyadic human interactions,” in *Acoustics, Speech and Signal Processing (ICASSP), IEEE International Conference on*, 2012.
- [31] E. Mower, M. J. Matarić, and S. Narayanan, “A framework for automatic human emotion classification using emotion profiles,” *Audio, Speech, and Language Processing, IEEE Transactions on*, vol. 19, no. 5, 2011.
- [32] P. Rao and L. L. Kupper, “Ties in paired-comparison experiments: A generalization of the bradley-terry model,” *Journal of the American Statistical Association*, vol. 62, no. 317, 1967.
- [33] P. Donmez and J. G. Carbonell, “Optimizing estimated loss reduction for active sampling in rank learning,” in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008.
- [34] C. Gentile and M. K. Warmuth, “Linear hinge loss and average margin,” in *NIPS*, 1998, vol. 11.
- [35] C. Burges, T. Shaked, E. Renshaw, A. Lazier, M. Deeds, N. Hamilton, and G. Hullender, “Learning to rank using gradient descent,” in *Proceedings of the 22nd international conference on Machine learning*. ACM, 2005.
- [36] J. D. Rennie and N. Srebro, “Loss functions for preference levels: Regression with discrete ordered labels,” in *Proceedings of the IJCAI multidisciplinary workshop on advances in preference handling*. Kluwer Norwell, MA, 2005.
- [37] T. Hastie, R. Tibshirani, et al., “Classification by pairwise coupling,” *The annals of statistics*, vol. 26, no. 2, 1998.

- [38] J. A. Hartigan and M. A. Wong, “Algorithm as 136: A k-means clustering algorithm,” *Applied statistics*, 1979.
- [39] P. Cortez, A. Cerdeira, F. Almeida, T. Matos, and J. Reis, “Modeling wine preferences by data mining from physicochemical properties,” *Decision Support Systems*, vol. 47, 2009.
- [40] A. Asuncion and D. Newman, “UCI machine learning repository,” 2007.
- [41] F. Alimoglu, D. Doc, E. Alpaydin, and Y. Denizhan, “Combining multiple classifiers for pen-based handwritten digit recognition,” 1996.
- [42] A. Trajman and R. Luiz, “McNemar  $\chi^2$  test revisited: comparing sensitivity and specificity of diagnostic examinations,” *Scandinavian journal of clinical and laboratory investigation*, vol. 68, no. 1, pp. 77–80, 2008.
- [43] D. Bone, M. P. Black, A. Ramakrishna, R. Grossman, and S. S. Narayanan, “Acoustic-prosodic correlates of awkward prosody in story retellings from adolescents with autism,” in *Proceedings of Interspeech*, Sept. 2015.
- [44] R. B. Grossman, L. R. Edelson, and H. Tager-Flusberg, “Emotional facial and vocal expressions during story retelling by children and adolescents with high-functioning autism,” *Journal of Speech, Language, and Hearing Research*, vol. 56, 2013.
- [45] D. Koller and N. Friedman, *Probabilistic graphical models: principles and techniques*, MIT press, 2009.

## Appendix 1: Proof for equation (10)

To prove:

$$q(z_{ij}^*) = p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) =$$

$$\left( p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \times \prod_{k=1}^K p(z_{ij}^k | z_{ij}^*, r_k) \right) / p(z_{ij}^1, \dots, z_{ij}^K)$$

Proof:

$$\begin{aligned} q(z_{ij}^*) &= p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) \\ &= p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) / p(z_{ij}^1, \dots, z_{ij}^K) \end{aligned} \quad (26)$$

By Bayes theorem:

$$\begin{aligned} & p(z_{ij}^*, z_{ij}^1, \dots, z_{ij}^K | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) / p(z_{ij}^1, \dots, z_{ij}^K) \\ &= p(z_{ij}^1, \dots, z_{ij}^K | z_{ij}^*, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) \\ &\times p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) / p(z_{ij}^1, \dots, z_{ij}^K) \end{aligned} \quad (27)$$

Based on the graphical model in Figure 1(c), we can say that  $z_{ij}^1, \dots, z_{ij}^K$  are independent of the attributes  $\mathbf{x}_i, \mathbf{x}_j$  and SVR vector  $\mathbf{w}$ , using the “indirect evidential effect” clause in [45].

$$\begin{aligned} & p(z_{ij}^1, \dots, z_{ij}^K | z_{ij}^*, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) = \\ & p(z_{ij}^1, \dots, z_{ij}^K | z_{ij}^*, r_1, \dots, r_K) \end{aligned} \quad (28)$$

Next, applying the “common clause” effect [45] to the graphical model in Figure 1(c), we can say that  $z_{ij}^1, \dots, z_{ij}^K$  are mutually independent given  $z_{ij}^*$ . Consequentially,  $z_{ij}^k$  is also independent of all  $r_{k'}$  for all  $k' \neq k$  due to the “common clause” effect. Therefore:

$$p(z_{ij}^1, \dots, z_{ij}^K | z_{ij}^*, r_1, \dots, r_K) = \prod_{k=1}^K p(z_{ij}^k | z_{ij}^*, r_k) \quad (29)$$

We can also say that  $z_{ij}^*$  is independent of  $r_1, \dots, r_K$  when the probability distribution is not conditioned on  $z_{ij}^1, \dots, z_{ij}^K$  again using the “common clause” effect [45].

$$p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) = p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \quad (30)$$



Replacing (29) and (30) into (27), we obtain

$$q(z_{ij}^*) = p(z_{ij}^* | z_{ij}^1, \dots, z_{ij}^K, \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) \quad (31)$$

$$= p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \prod_{k=1}^K p(z_{ij}^k | z_{ij}^*, r_k) / p(z_{ij}^1, \dots, z_{ij}^K) \quad (32)$$

## Appendix 2: Proof for equation (16)

To prove:

$$\begin{aligned} & \log p(z_{ij}^*, z_{ij}^1, \dots, z_{ij}^K | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}, r_1, \dots, r_K) \\ &= \sum_{k=1}^K \log p(z_{ij}^k | z_{ij}^*, r_k) + \log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \end{aligned} \quad (33)$$

*Proof:*

Application of (27)-(30) to the left hand side of (33) yields the desired result.

## Appendix 3: Probability distribution for optimization in equation (18)

The goal in the M-step of the EM algorithm in order to obtain  $\mathbf{w}$  was to perform the following optimization.

$$\begin{aligned} \mathbf{w} &= \arg \max_{\mathbf{w}} \mathcal{M}_{\mathbf{w}} \\ &= \arg \max_{\mathbf{w}} \sum_{z_{ij}^* \in \{0,1\}} q(z_{ij}^*) \log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \end{aligned} \quad (34)$$

Where

$$p(z_{ij}^* = 1 | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) = \frac{\exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle}{1 + \exp \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle} \quad (35)$$

$$p(z_{ij}^* = 0 | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) = 1 - p(z_{ij}^* = 1 | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}) \quad (36)$$

Instead, we performed the optimization in (18), restated below.

$$\begin{aligned} \mathbf{w} &= \arg \min_{\mathbf{w}} \mathcal{M}_{\mathbf{w}} = \arg \min_{\mathbf{w}} \left( q(z_{ij}^* = 1) [1 - \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle]_+ \right. \\ &\quad \left. + q(z_{ij}^* = 0) [1 - \langle \mathbf{w}, \{\mathbf{x}_j - \mathbf{x}_i\} \rangle]_+ \right) \end{aligned} \quad (37)$$

Above optimization can be rewritten as shown in (38).

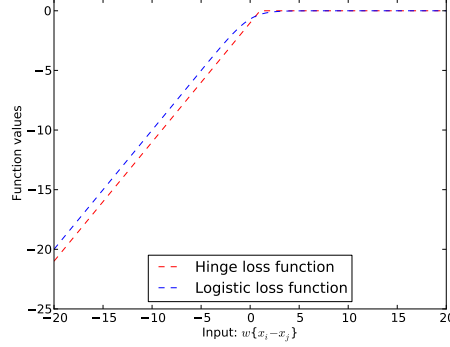


Figure 4: Plot comparing the values of the negative hinge loss function  $(-1 \times [1 - \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle]_+)$  and the log of logistic loss function  $(\log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}))$ .

$$\begin{aligned} \mathbf{w} = \arg \max_{\mathbf{w}} & \left( q(z_{ij}^* = 1) (-1 \times [1 - \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle]_+) \right. \\ & \left. + q(z_{ij}^* = 0) (-1 \times [1 - \langle \mathbf{w}, \{\mathbf{x}_j - \mathbf{x}_i\} \rangle]_+) \right) \end{aligned} \quad (38)$$

We compare the negative hinge loss function  $(-1 \times [1 - \langle \mathbf{w}, \{\mathbf{x}_i - \mathbf{x}_j\} \rangle]_+)$  and the log of the logistic loss function  $(\log p(z_{ij}^* | \mathbf{x}_i, \mathbf{x}_j, \mathbf{w}))$  stated in (34). Figure 4 shows the values that these functions take with respect to the input  $\mathbf{w}\{\mathbf{x}_i - \mathbf{x}_j\}$ . The plots indicate that the values taken by the two functions are very close to each other. One difference is around an input value of 0, where the hinge loss function is not differentiable but the logistic loss function is. More importantly, the slopes of the two functions are same for a large range of input and therefore, for all practical purposes, the gradient descent algorithm should provide similar results after replacing the logistic loss function with hinge loss function in the M-step of the EM algorithm. However, we were unable to theoretically prove that the algorithm still falls under the paradigm of generalized EM algorithm, and therefore is an approximation in the EM algorithm.